

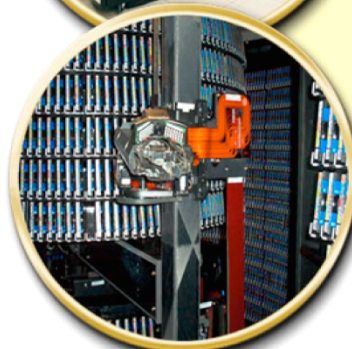
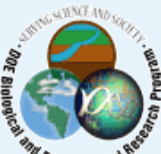
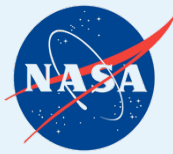
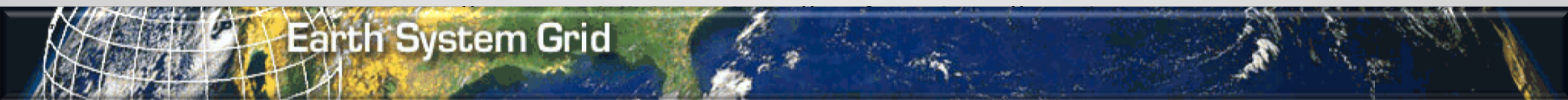
Partnership for 21st Century Earth System Science

Data Federation and Analysis Integration from Multiple Interagency Sources



Dean N. Williams

On behalf of the ESG-CET, ESGF, UV-CDAT, and Visual Data Explore
Teams



Goals

The World's Source for Climate Science Data

- Make data more useful to researchers and policy makers by developing collaborative technology that enhances data usability
- Meet the specific needs of national and international climate projects for distributed databases, data access, and data movement
- Provide a universal and secure Web-based data access portal for broad-based multi-model, observational, and reanalysis data collections
- Provide a wide range of climate data-analysis tools and diagnostic methods to international climate centers and U.S. government agencies.

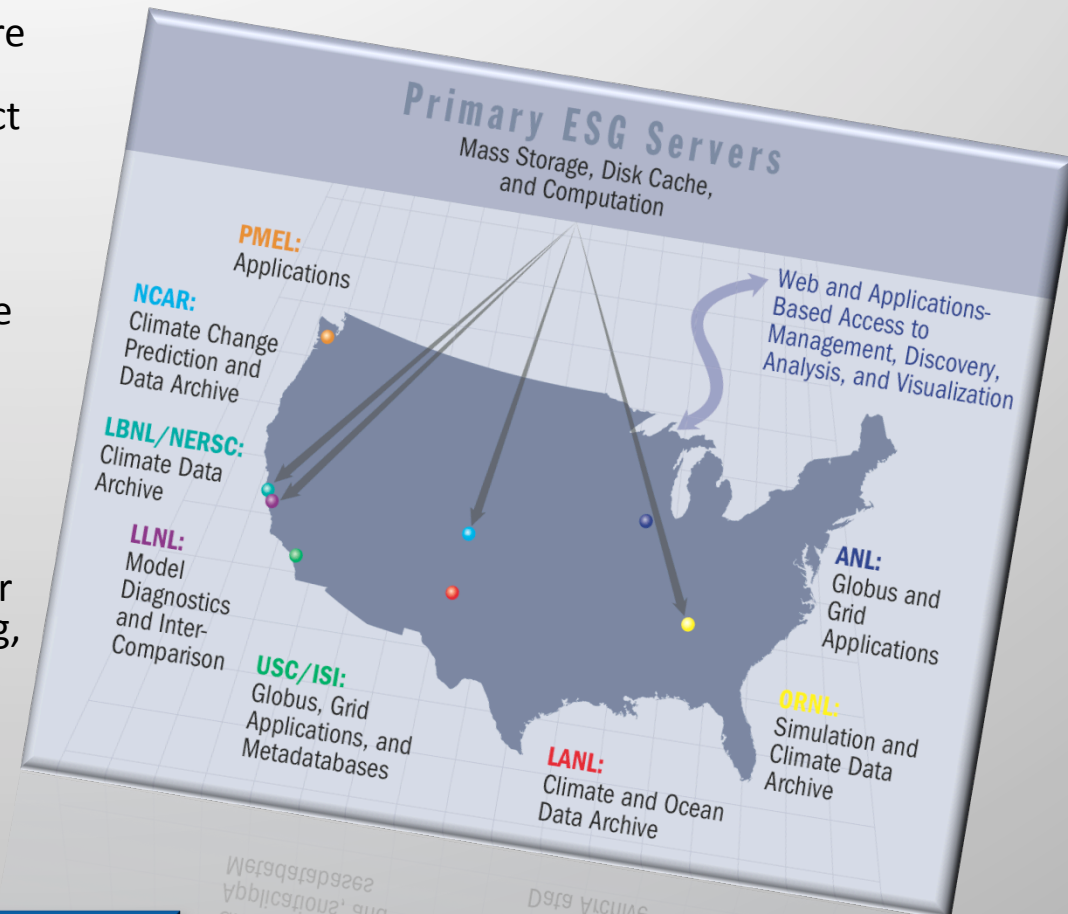


“Focused on broadening and strengthening the Earth System Science community by supporting a diverse set of national and international data collections.”

The Earth System Grid Federation (ESGF)

Building a Global Infrastructure for Climate Change Research

- ESGF is a free, open consortium of institutions, laboratories and centers around the world that are dedicated to supporting research of Climate Change, and its environmental and societal impact
- Historically originated from Earth System Grid (ESG) project, expanded beyond its constituency and mission to include many other partners in the U.S., Europe, Asia, and Australia
- Global Organization for Earth System Science Portal (GO-ESSP)** Groups working at many projects: ESG, Earth System Curator, Metafor, Global Interoperability Program, Infrastructure for the European Network for Earth System Modeling, and many more
- U.S. funding from DOE, NASA, NOAA, NSF

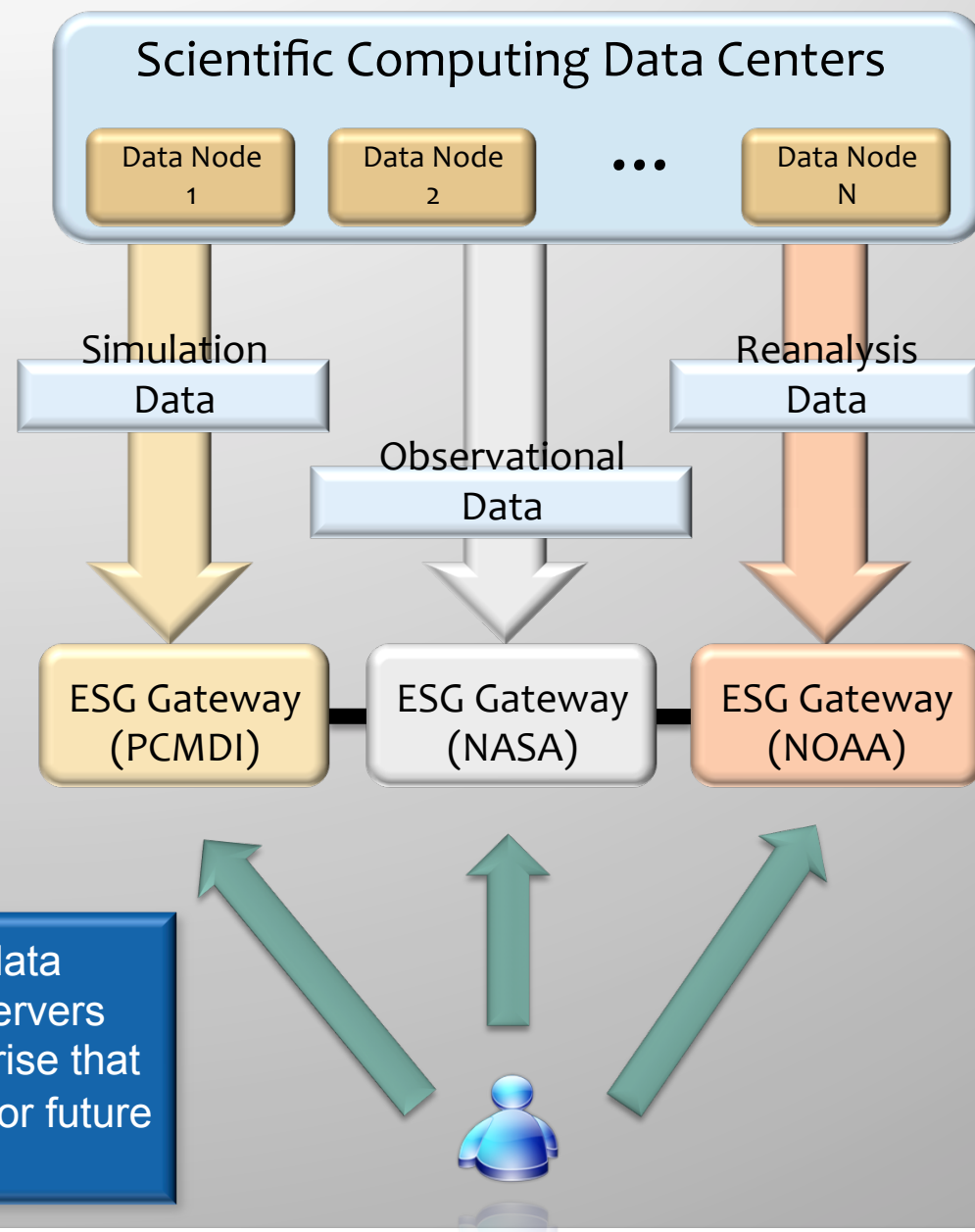


“Providing climate researchers worldwide with access to: data, information, models, analysis tools, and computational resources required to make sense of enormous Earth System data sets”

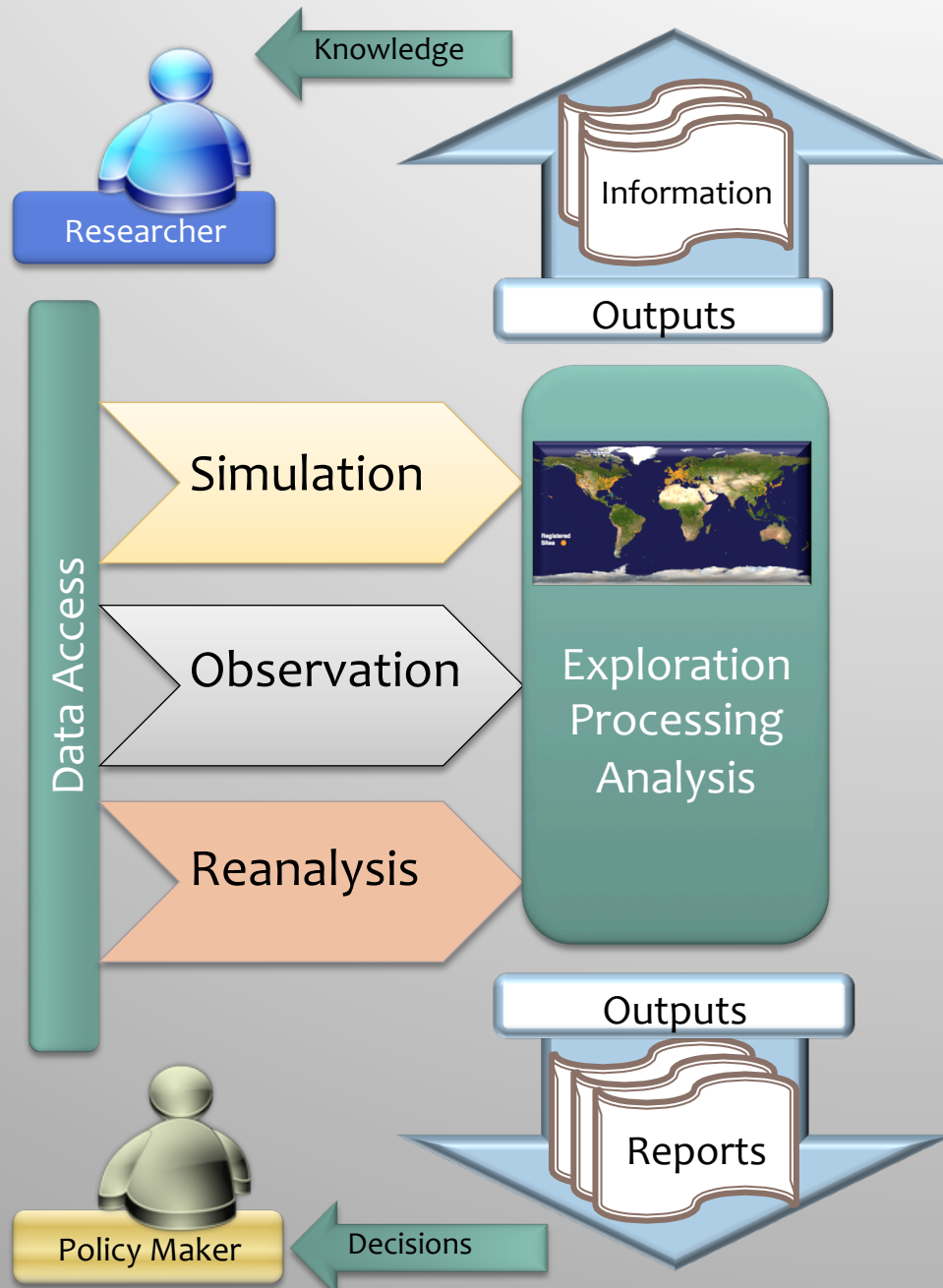
ESGF diverse partners and data sets

The World's Source for Earth System Science Data

- U.S.: ANL, ESRL, GFDL, NASA, NOAA, LANL, LBNL, PCMDI/LLNL, NCAR, ORNL, PMEL, USC/ISI, RPI
- Europe: BADC, UK-MetOffice, DKRZ, MPIM, IPSL, LSCE
- Asia: Univ. of Tokyo, Japanese Centre for Global Environmental Research, Jamstec, Korea Meteorological Administration
- Australia: ANU, Australian Research Collaboration Service, Government Department of Climate Change, Victorian Partnership for Advanced Computing, Australian Environment and Resource Management
- .. and many more ...



The Climate Community is making revolutionary end to end changes in Data Integration



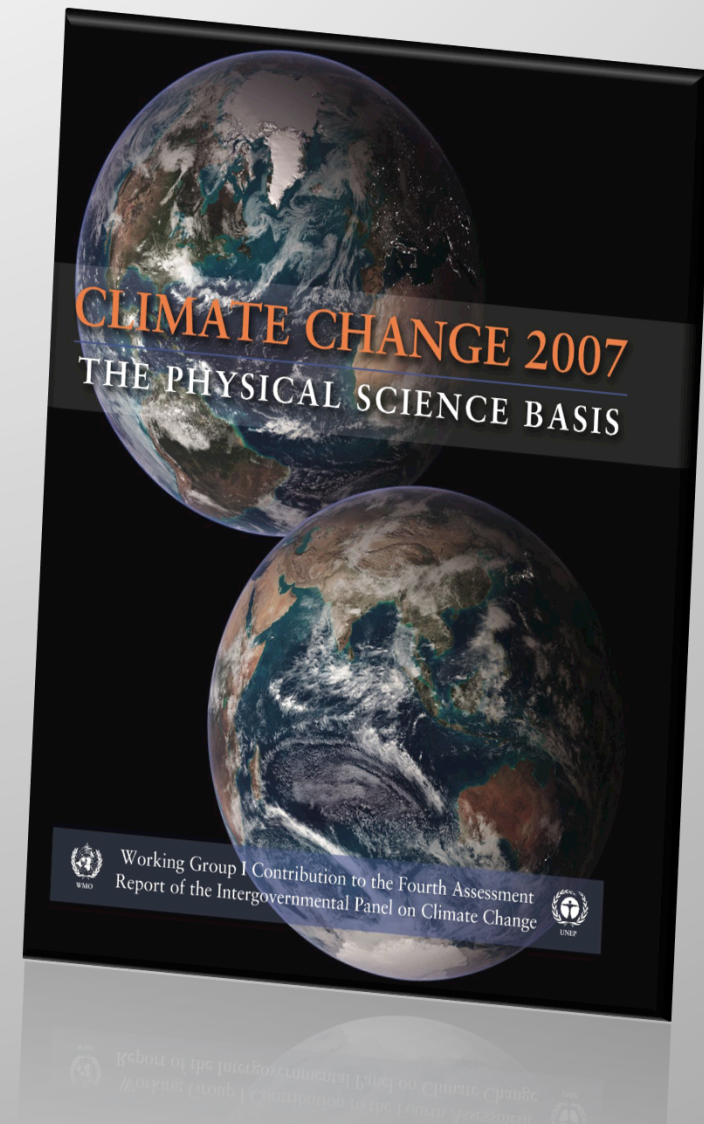
The challenge is integration:

- Massive data archives (many PB moving to XB)
- Multiple data centers worldwide
 - Existing IT infrastructure and separate security domains
- Heterogeneous data sources (models, observations, reanalysis)
- Multiple physical realms (atmosphere, ocean, land, sea ice)
- Multiple data, metadata formats and conventions
- Multiple scales (global, regional and local)
- Cyber security
- Multiple audiences (scientists, policy makers, students, educators)

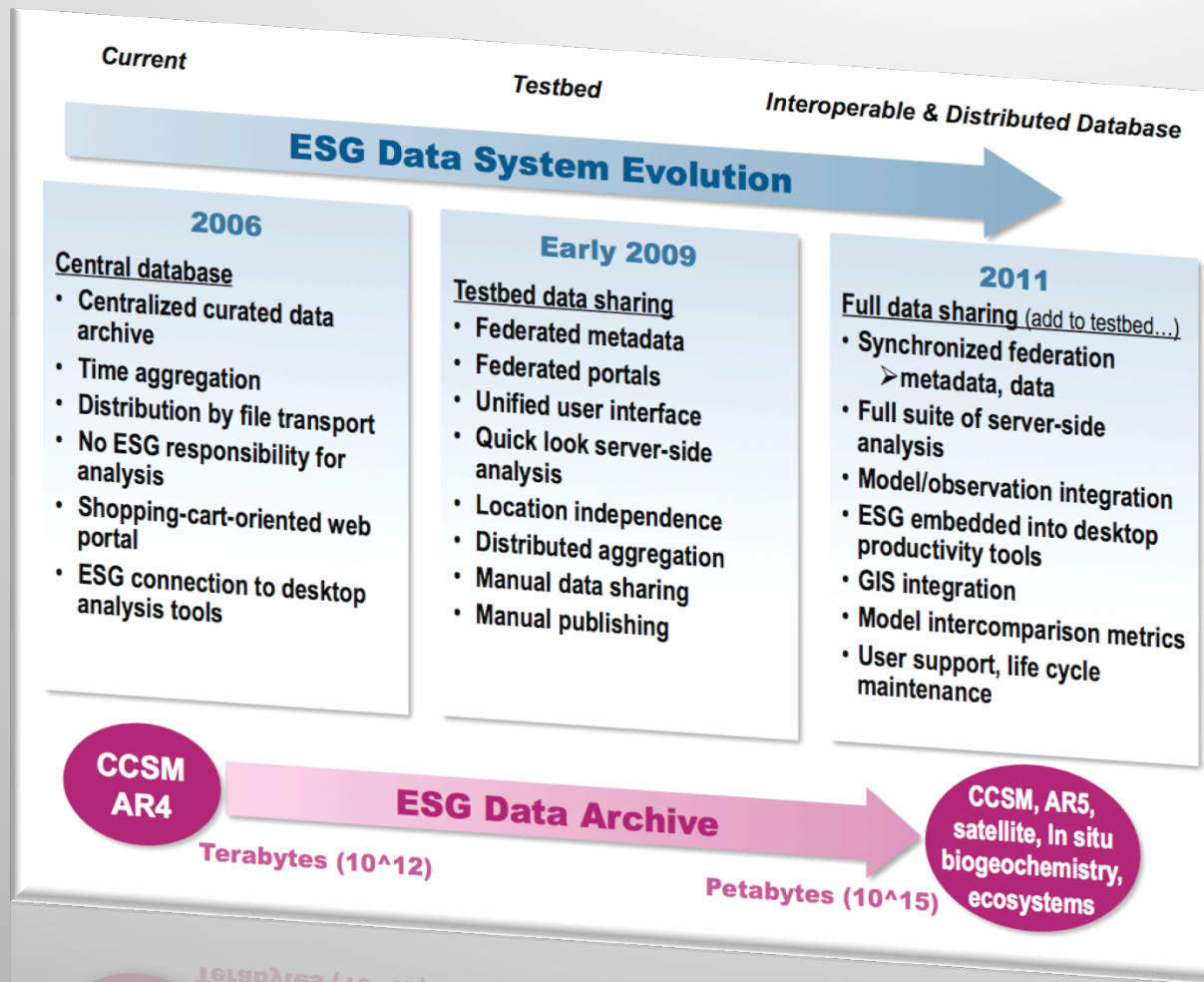
ESG is positioned to wisely build capabilities to meet climate change, policy, and environmental data needs

- Broad investments in climate change research
 - Development of climate models
 - Climate change simulation
 - Observational programs
 - Reanalysis archives
 - Model intercomparisons
 - Uncertainty Quantification and Testbeds
- Climate change research is increasingly data-intensive
 - Analysis and intercomparison of simulation and observations from many sources
 - Data used by model developers, policy makers, health officials, etc.
 - Working group I (researchers)
 - Working group II (impacts, adaptation, and vulnerability)
 - Working group III (mitigation of climate change)
- Broad Impact of ESG
 - Over 25K users

“For leadership in implementing, maintaining, and facilitating access to the CMIP3 multi-model data set archive, which led to a new era in climate system analysis and understanding.” – AMS Award 2010



Coupled Model Intercomparison Project, Phase 5 (CMIP5) – Major driver behind ESGF development



- CMIP5 multi-model archive expected to include
 - 3 suites of experiments
 - 40+ models
 - ~10 PB distributed data
 - ~2 PB central replicated archive
 - Contribution from 25+ modeling centers in 17+ countries
- Driver for scale of data, global distribution
- Timeline fixed by IPCC (2011 – 2013)
- Working with over 30 key national and international partners to establish ESGF

“International climate modeling activities are largely coordinated by the WCRP/CLIVAR Working Group on Coupled Modeling (WGCM). ”

CMIP3 archive vs. CMIP5 archive

Modeling group		CMIP3 volume (GB)
NCAR	USA	9,172.8
MIROC3	Japan	3,974.9
GFDL	USA	3,842.5
IAP	China	2,867.7
MPI	Germany	2,699.5
CSIRO	Australia	2,088.2
CCCMA	Canada	2,070.6
INGV	Italy	1,472.2
GISS	USA	1,096.8
MRI	Japan	1,024.5
CNRM	France	999.1
IPSL	France	997.7
UKMO	UK	972.8
BCCR	Norway	861.9
MIUB	Germany/Korea	477.2
INMCM3	Russia	368.2
Totals		34,986.6

Modeling group		CMIP5 volume (GB)
MPI	Germany	710,000
NCAR	USA	410,000
MRI	Japan	312,000
GFDL	USA	151,000
MIROC3	Japan	115,000
UKMO	UK	89,000
CNRM	France	64,000
IAP	China	63,000
U Reading	UK	63,000
EC	Europe	50,000
GISS	USA	50,000
INGV	Italy	50,000
IPSL	France	45,000
INMCM3	Russia	32,000
NorClim	Norway	30,000
CCCMA	Canada	29,000
CAWCR	Australia	21,000
CSIRO	Australia	20,000
METRI	Korea	13,000
Totals		2,317,000

The ESGF system is a network of Gateways and Data Nodes

■ Gateways

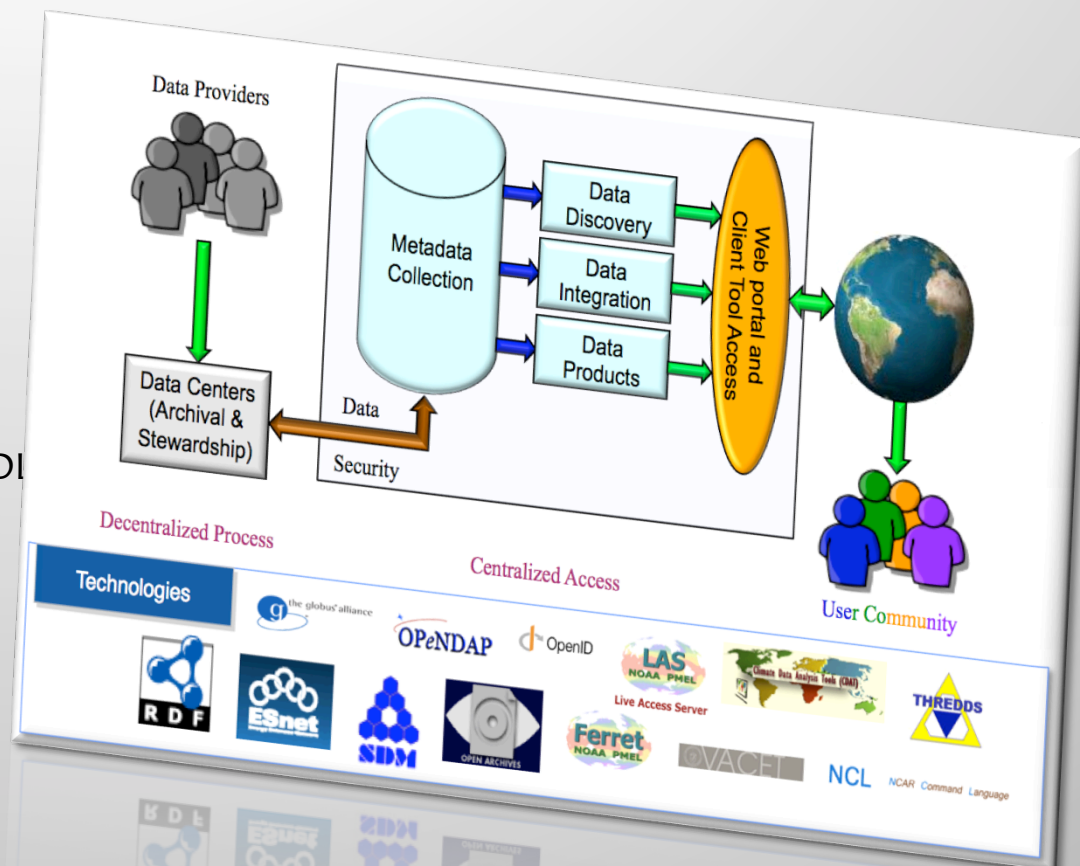
- Where data is discovered, requested
- Portals, search capability, distributed metadata, registration and user management
- May be customized to an institution's requirements, topical focus
- Fewer sites than Data Nodes
- Currently: PCMDI, NCAR, ORNL, LBNL, NASA, BADC, DKRZ, ANU; coming soon: ANL, PNNL, GFDL, JAMSTEC

■ Data Nodes

- Where data are stored and published
- Data may be on disk or tertiary mass store
- Each Data Node can publish to any Gateway (facilitates topical Gateways)
- Data reduction/analysis
- Complex architecture, including possible minimalist deployment without services
- Anticipates ~20 Data Nodes for CMIP5, many others have expressed interest (over 50 sites)

■ Sites

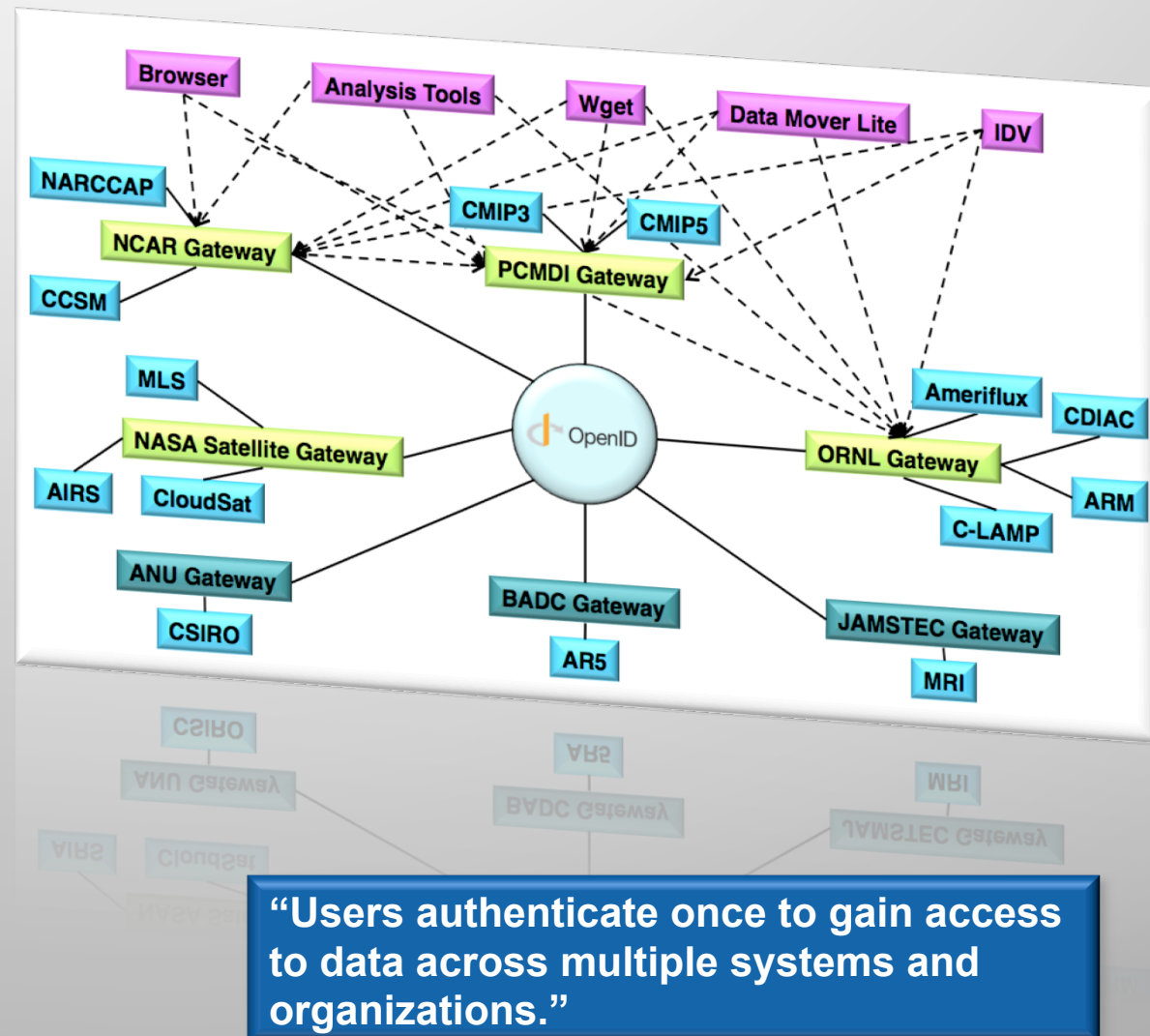
- A site can be both a Gateway and a Data Node



“Multi-disciplinary teams are required to attack large-scale problems: ESG, Earth System Curator, Metafor, Global Interoperability Program, Infrastructure for the European Network for Earth System Modeling, and many more.”

ESGF is a distributed data archival and retrieval system

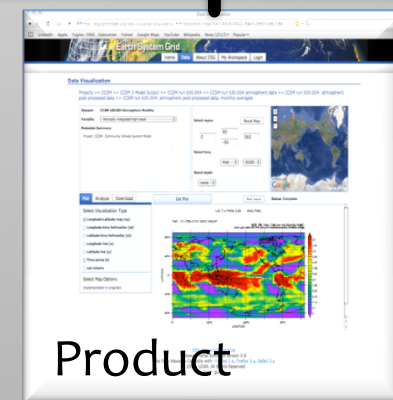
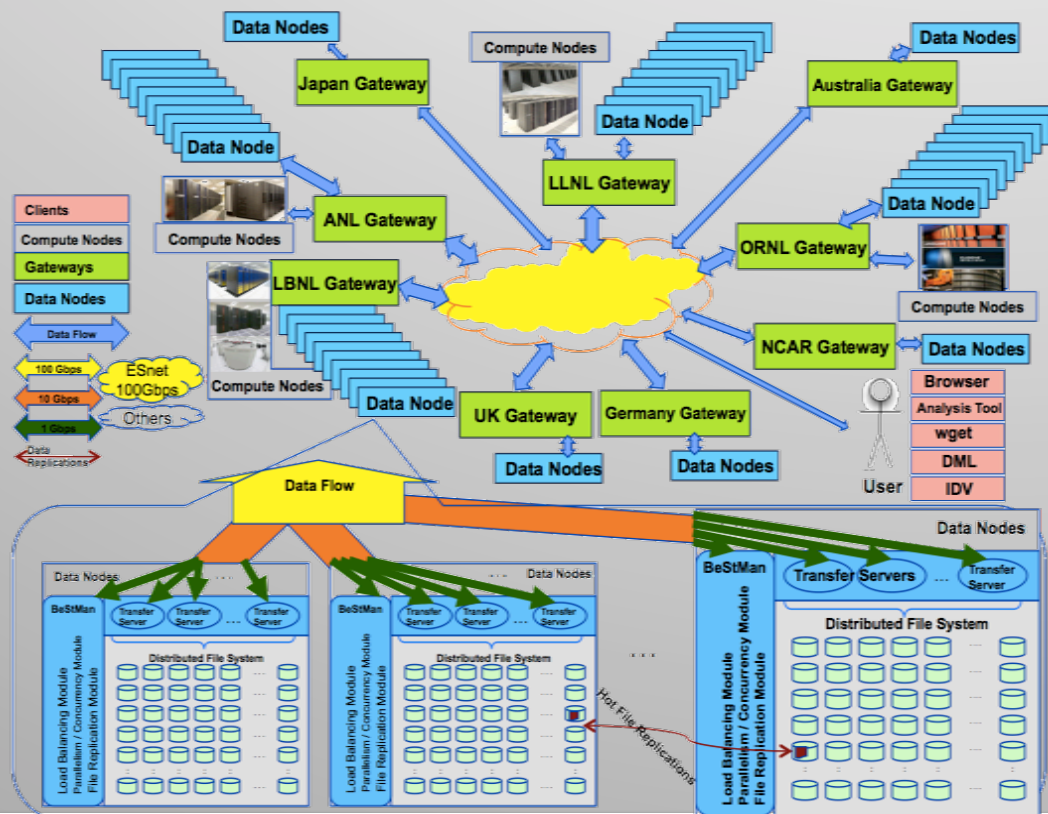
- Distributed and federated architecture
- Support discipline specific Gateways
- Support browser-based and direct client access
- Single Sign-on
- Automated GUI-based publication tools
- Full support for data aggregations
 - A collection of files, usually ordered by simulation time, that can be treated as a single file for purposes of data access, computation, and visualization
- User notification service
 - Users can choose to be notified when a data set has been modified



ESG operations– Improved users access and server-side analysis capabilities

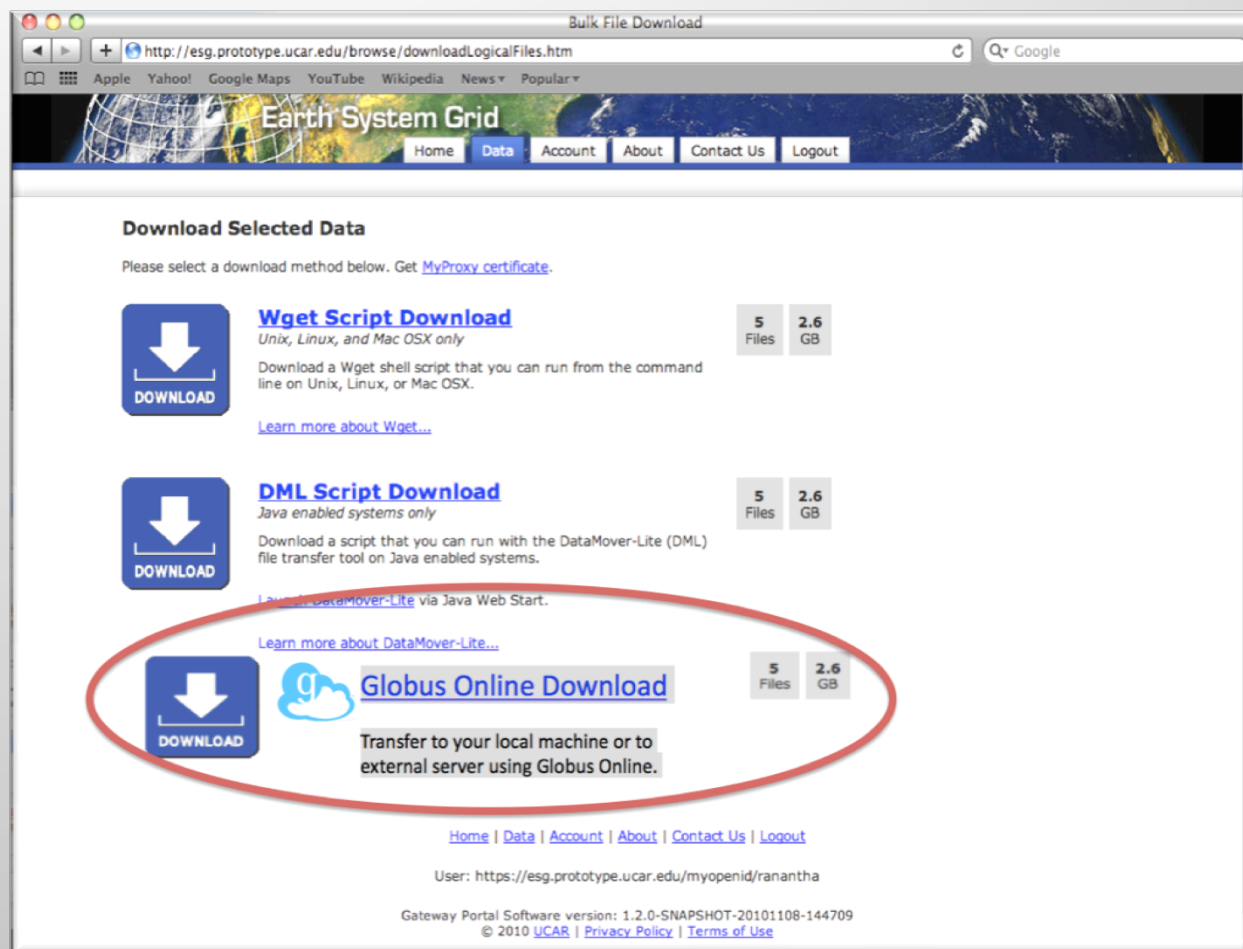
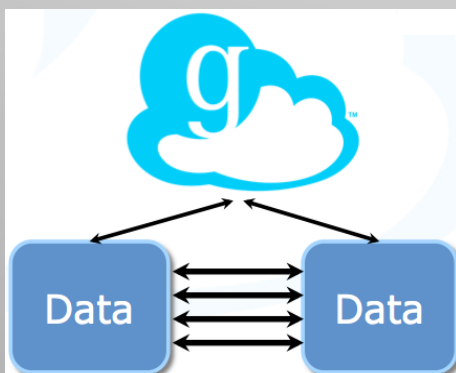
Gateway URLs

- <http://pcmdi3.llnl.gov/esgcet/home.htm>
- <http://www.earthsystemgrid.org>
- <http://esg.ccs.ornl.gov/esgcet/home.htm>
- <http://esg-gateway.jpl.nasa.gov/home.htm>



Download selected data

- Easy “fire and forget” file transfers
- Automatic fault recovery
- High performance
- Simplify use of multiple security domains
- No client software installation
- New features automatically available
- Consolidated support and troubleshooting



Gateway data discovery

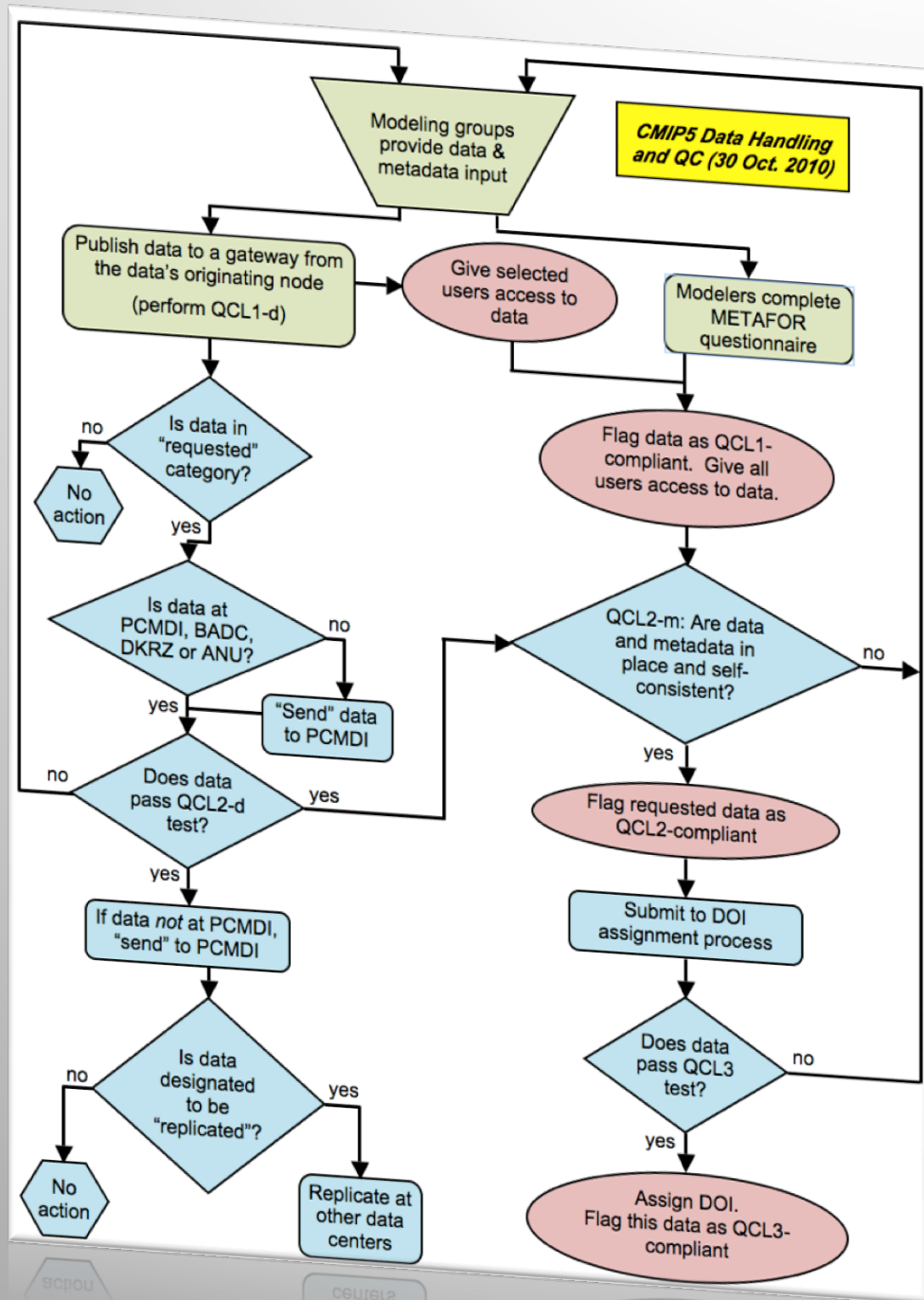
The operations performed within the ESG infrastructure

- Faceted search – allows users to easily and rapidly browse data sets
- Data Reference Syntax (DRS) – identify data sets wherever they might be located within the distributed ESGF archive
- Data versioning – the tracking of multiple versions of data
- Data migration – the procedure of migrating the relational database content when upgrading the underlying schema (Liquibase)
- Attribute and Authorization service - Security Assertion Markup Language (SAML)-signed assertions in response to attribute and authorization queries by remote clients
- Metadata exchange - based on the Open Archives Initiative's Protocol for Metadata Harvesting (OAI-PMH), has been updated to support versioning and replication, as well as the capability to execute selected harvesting by project
- Data download - generation of wget scripts and integration with the DataMover-Lite (DML) desktop client and Globus Online
- Model metadata - Earth System Curator project has continued to work with ESG to expand the gateway functionality for ingesting and servicing model metadata , including full handling of the CMIP5 conformance properties, ingestion of Common Information Model (CIM) metadata from the Common Metadata for Climate Modelling Digital Repositories (METAFOR) Questionnaire application
- Federated authentication - federated system allowing user access via the ESG gateways and supporting interoperability with other non-ESG partner data centers
- Gateway user interfaces - pages for metadata search, data download, model traceback, user and group administration

The operations performed by the data node within the ESG infrastructure

- [illegible]

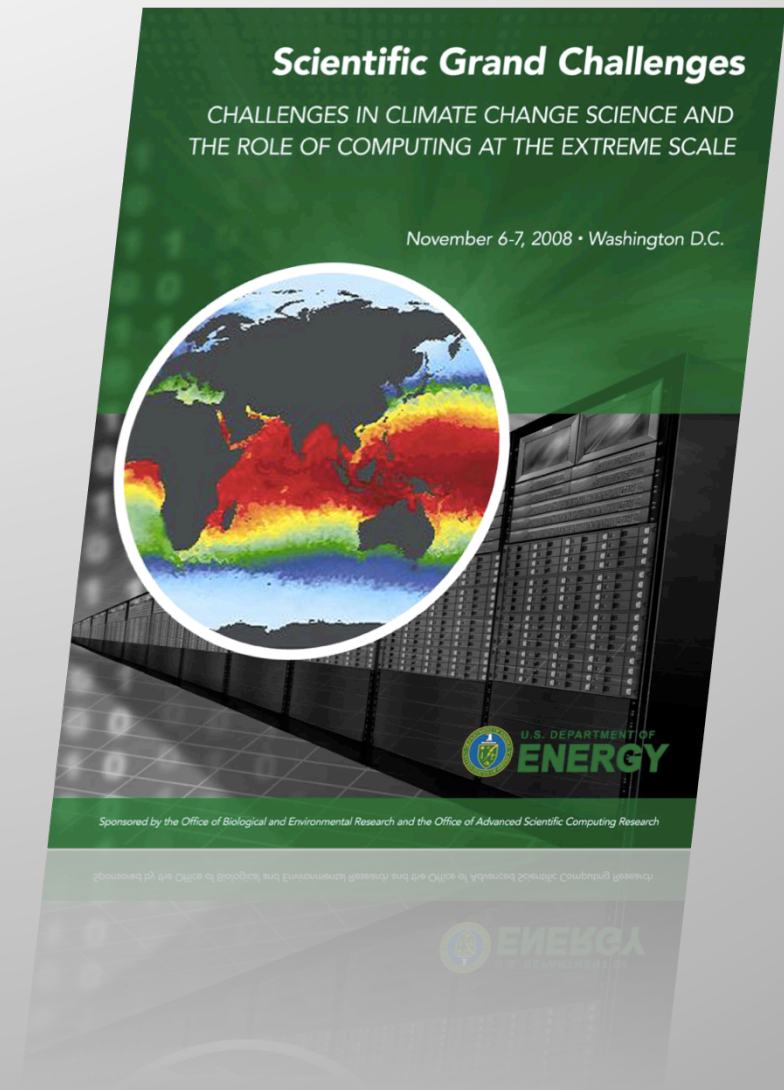
CMIP5 ESG Data flow operations



- Publishing data to an ESG Gateway performs QC Level 1 (QCL1) check
 - QCL1 data are visible to users and are identified as QCL1 on the UI
- DKRZ (MPI) quality control code is run on data to perform QC Level 2 (QCL2) check
 - QCL2 data are visible to users and are identified as QCL2 on the UI
- Visual inspections are performed for inconsistencies and metadata correctness at QC Level 3 (QCL3) check
 - QCL3 data are visible to users and are identified as QCL3 on the UI
 - Digital Object Identifiers (DOIs) are given to data sets that pass the QCL3 check

Critical Functions - Enhancing data uses, measurement, and validation

- Replication – A node can choose to **replicate a collection of the data sets** published by a different node. This includes replicas of aggregated data sets.
- Metrics – **Quantitative measurement** and reporting of the usage and performance of the ESG Enterprise system.
- Bulk Data Movement (BDM) – The **transfer of large collections** between sites reliably and with good performance and with a high-level of easy-of-use.
- **Provenance** - important for science because it **helps to interpret and reproduce the results of an experiment**; to understand the chain of reasoning used in the production of a result; to verify that the experiment was performed according to acceptable procedures, to track who performed the experiment and who is responsible for its results.
- **Product Service via the Live Access Server (LAS)** – Access to analysis & visualization tools from distributed sources in an secure environment.
- **Easy Software Stack Installation** - Installation process is an interactive script that will fully install and configure all components. Also on the horizon are Virtual Machine installation (Linux: Centos, SUSE)

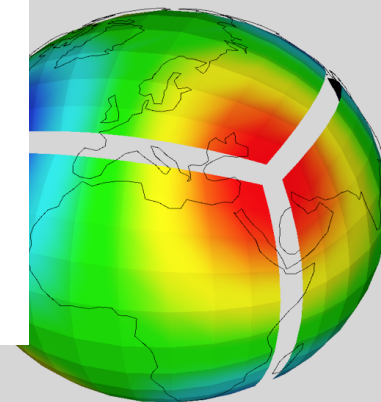
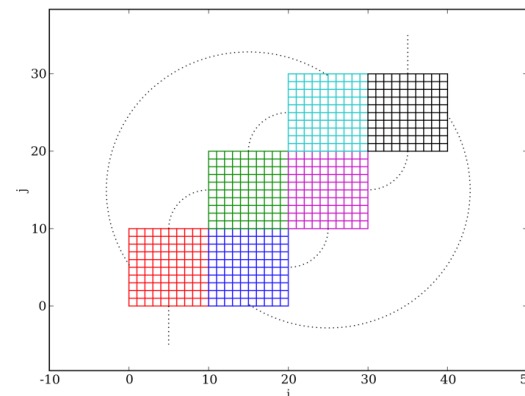


Metadata Requirements revolve around the netCDF Climate and Forecast (CF) Convention

- Each File contains only a single output field
- Most data are a function of longitude and latitude, represented as a Cartesian product of longitude and latitude axes
 - Gridspec – a standard for description of grids used in Earth System models; Gridspec is designed for inclusion within the CF metadata convention.
 - LibCF – supports the creation of scientific data files conforming to the CF conventions, using the netCDF API.
 - UV-CDAT is one of several standard analysis tools that are needed to operate on the Earth System Grid's (ESG's) "Product Service" back-end. This tool (along with ESG) conforms to the CF convention.

Table A1a: Monthly-mean 2-d atmosphere or land surface data (longitude, latitude, time:month).

	CF standard_name	output variable name	units	notes
1	air_pressure_at_sea_level	psl	Pa	
8	surface_snow_thickness	snd	m	this thickness when multiplied by the average area of the grid cell covered by snow yields the time-mean snow volume. Thus, for time means, compute as the weighted sum of thickness (averaged over the snow-covered portion of the grid cell) divided by the sum of the weights, with the weights equal to the area covered by snow. report as 0.0 in snow-free regions.
15	surface_temperature	ts	K	"skin" temperature (i.e., SST for open ocean)
16	surface_air_pressure	ps	Pa	not mean sea-level pressure
19	atmosphere_water_vapor_content	prw	kg m ⁻²	vertically integrated through the atmospheric column
21	surface_runoff_flux	mrros	kg m ⁻² s ⁻¹	compute as the total surface runoff leaving the land portion of the grid cell divided by the land area in the grid cell; report as "missing" or 0.0 where the land fraction is 0.
22	runoff_flux	mrro	kg m ⁻² s ⁻¹	compute as the total runoff (including "drainage" through the base of the soil model) leaving the land portion of the grid cell divided by the land area in the grid cell; report as "missing" or 0.0 where the land fraction is 0.



MoDAVE

3 files of cubed sphere

Climate Model Output Rewriter (CMOR) – Providing computing and data for CA power

- A C library (with FORTRAN 90 bindings) for rewriting model output and ensure compliance with the IPCC requirements.
- Relies on an input file to supply much of the “standard metadata associated with the IPCC standard output fields.
- The CMOR input files can be reconfigured to change the requirements and meet the needs of other model intercomparison projects. Capabilities to:
 - Reorder dimensions;
 - Reverse the order of coordinate values;
 - Automatic units conversion and appropriate data scaling; etc
 - Checker
- Data Reference Syntax (DRS)
 - The common naming system should also be used in files, directories, metadata and URLs to facilitate standard identification of data sets.
- URLs:
 - <http://cf-pcmdi.llnl.gov>
 - <http://www2-pcmdi.llnl.gov/cmor>
 - <http://cmip-pcmdi.llnl.gov>
 - <http://cmip-pcmdi.llnl.gov/cmip5/docs/>

▪ Required attributes for variables:

→ associated_files = a string listing the base URL for CMIP5 and the location of the model's gridspec file, followed, as appropriate, by the name of the file containing the grid cell areas and/or grid cell volumes. For CMIP5 this string is: "baseURL:<http://cmip-pcmdi.llnl.gov/> <<http://www-pcmdi.llnl.gov/>> CMIP5/dataLocation gridspecFile:<gridspec file name> [areacella:<atmos. cell area file name>] [areacello:<ocean cell area file name>] [volcello:<ocean cell volume file name>]", where cell area and cell volume are only sometimes required. Here is an example:

→ associated_files=

→ "baseURL: <http://cmip-pcmdi.llnl.gov/> <<http://www-pcmdi.llnl.gov/>> CMIP5/
dataLocation

→ gridspecFile: gridspec_ocean_fx_IPSL-CM5_historical_r0i0p0.nc

→ areacello: areacello_fx_IPSL-CM5_historical_r0i0p0.nc

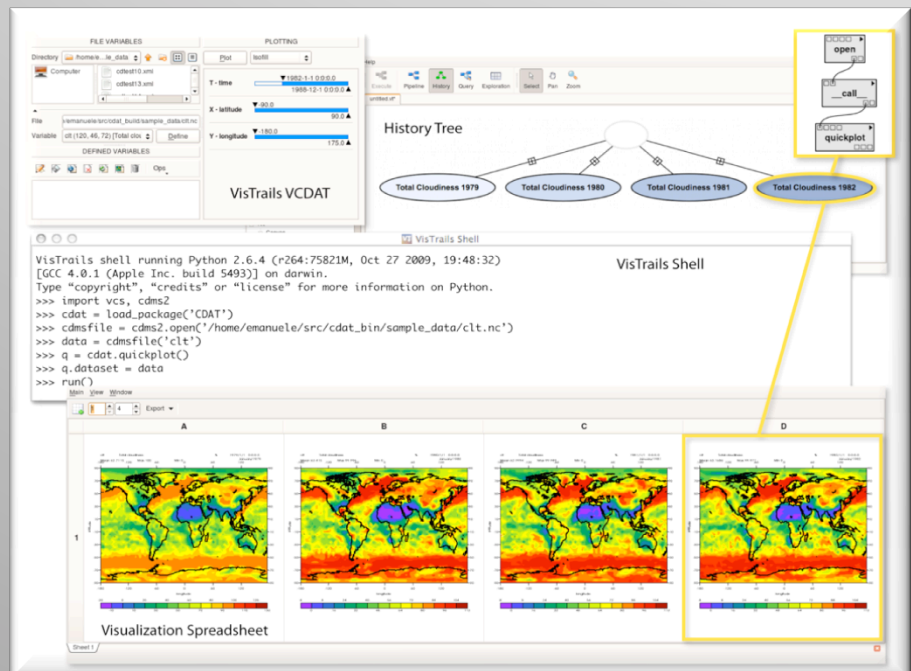
→ volcello: volcello_fx_IPSL-CM5_historical_r0i0p0.nc"

→ Note that the “cell_measures” column (U) of the CMIP5 Requested Output <http://pcmdi-cmip.llnl.gov/cmip5/docs/standard_output.pdf> tables indicates whether or not areacella, areacello, and/or volcello should be included as associated_files. For each model version and experiment, only one area field is requested by CMIP5 for the atmosphere and only one area and one volume field are requested for the ocean. These cell areas should be defined such that exact global integrals of energy fluxes at the surface and “top of the atmosphere” can be computed. It is understood that for some staggered grids, the meshes for horizontal velocities might be offset from the radiation points, so in these cases exact global integrals of momentum would require areas not requested by CMIP5. In the associated_files attribute, include references to areacella, areacello, and volcello only for variables carried on the same mesh as the areas and volumes (i.e., only when it is appropriate to do so).

More focused and leveraged partnerships

UV-CDAT

- Integrate DOE's climate modeling and measurements archives
- Develop infrastructure for national and international model/data comparisons
- Deploy a wide-range of climate data visualization, diagnostic, and analysis tools with familiar interfaces for very large, high resolution climate data sets
- Workflow – data flows are directed graphs describing computational tasks
- Takes advantage of ESG data management



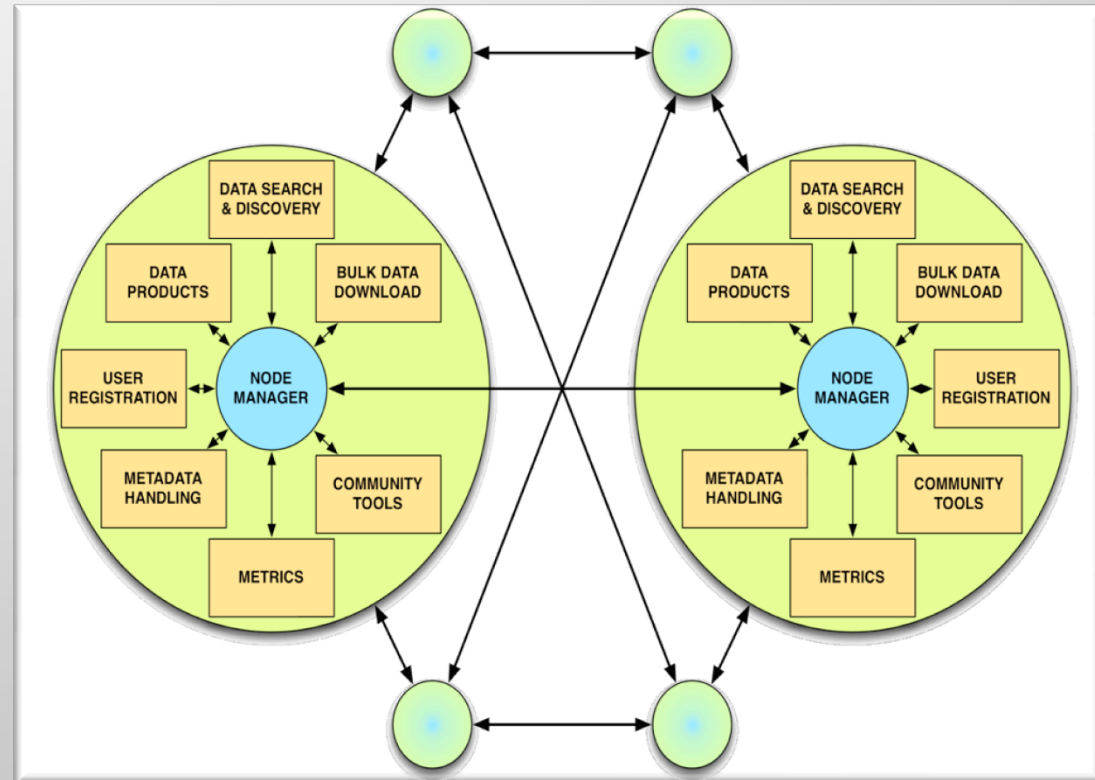
Visual Data Explorer

- Powerful visualization techniques with climate data sets involving simulation, observational, and reanalysis data
- Parallel visualization applications
- Python scripting
- Takes advantage of remote ESG data management via UV-CDAT



Possible Future ESGF developments plans

- **Other scientific domains:** biology, chemistry, high-energy physics, etc.
- **Peer-to-Peer Architecture**
- **Analysis** services for:
 - extremely large data sets
 - multiple large data sets are not co-located
 - cloud computing
- **Data integration** and advanced metadata capabilities
- **Advanced product services** via multiple scripting languages
- Integration of security assertion Markup Languages (**SAML**) identity providers
- **Measuring replication and data access patterns** in extreme scale ESG
- **Workflow and provenance**
- **Virtual Organization** management as **Software as a Service**
- **Advanced networks** as easy-to-use community resources
- **Management of open source, community-driven software development**



Immediate next steps

- Build the CMIP5 (IPCC AR5) data repository
- Inclusion of observational and reanalysis data from NASA, NOAA, ORNL, others
- Community-driven open source software development and project governance

